



## Effective Fusion of Deep Multitasking Representations for Robust Visual Tracking

Zadeh, Seyed Mojtaba Marvasti; Ghanei-Yakhdan, Hossien; Kasaei, Shohreh ; Nasrollahi, Kamal; Moeslund, Thomas B.

*Published in:*  
Visual Computer

*DOI (link to publication from Publisher):*  
[10.1007/s00371-021-02304-1](https://doi.org/10.1007/s00371-021-02304-1)

*Publication date:*  
2022

*Document Version*  
Accepted author manuscript, peer reviewed version

[Link to publication from Aalborg University](#)

*Citation for published version (APA):*  
Zadeh, S. M. M., Ghanei-Yakhdan, H., Kasaei, S., Nasrollahi, K., & Moeslund, T. B. (2022). Effective Fusion of Deep Multitasking Representations for Robust Visual Tracking. *Visual Computer*, 38(12), 4397-4417. <https://doi.org/10.1007/s00371-021-02304-1>

### General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal -

### Take down policy

If you believe that this document breaches copyright please contact us at [vbn@aub.aau.dk](mailto:vbn@aub.aau.dk) providing details, and we will remove access to the work immediately and investigate your claim.

---

©2020 Springer. Personal use of this material is permitted. Permission from Springer must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

# Effective Fusion of Deep Multitasking Representations for Robust Visual Tracking

Seyed Mojtaba Marvasti-Zadeh ·  
Hossein Ghanei-Yakhdan · Shohreh Kasaei ·  
Kamal Nasrollahi · Thomas B. Moeslund

Received: date / Accepted: date

**Abstract** Visual object tracking remains an active research field in computer vision due to persisting challenges with various problem-specific factors in real-world scenes. Many existing tracking methods based on discriminative correlation filters (DCF) employ feature extraction networks (FENs) to model the target appearance during the learning process. However, using deep feature maps extracted from FENs based on different residual neural networks (ResNets) has not previously been investigated. This paper aims to evaluate the performance of twelve state-of-the-art ResNet-based FENs in a DCF-based framework to determine the best for visual tracking purposes. First, it ranks their best feature maps and explores the generalized adoption of the best ResNet-based FEN into another DCF-based method. Then, the proposed method extracts deep semantic information from a fully convolutional FEN and fuses it with the best ResNet-based feature maps to strengthen the target representation in the learning process of continuous convolution filters. Finally, it introduces a new and efficient semantic weighting method (using semantic segmentation feature maps on each video frame) to reduce the drift

---

• S. M. Marvasti-Zadeh

Digital Image and Video Processing Lab (DIVPL), Department of Electrical Engineering, Yazd University, Yazd, Iran. He is also a member of Image Processing Lab (IPL), Sharif University of Technology, Tehran, Iran, and Vision and Learning Lab, University of Alberta, Edmonton, Canada (E-mail: mojtaba.marvasti@ualberta.ca).

• H. Ghanei-Yakhdan (corresponding author)

Digital Image and Video Processing Lab (DIVPL), Department of Electrical Engineering, Yazd University, Yazd, Iran (E-mail: hghaneiy@yazd.ac.ir).

• S. Kasaei

Image Processing Lab (IPL), Department of Computer Engineering, Sharif University of Technology, Tehran, Iran (E-mail: kasaei@sharif.edu).

• K. Nasrollahi

Visual Analysis of People Lab (VAP), Department of Architecture, Design, and Media Technology, Aalborg University, Aalborg, Denmark (E-mail: kn@create.aau.dk).

• T. B. Moeslund

Visual Analysis of People Lab (VAP), Department of Architecture, Design, and Media Technology, Aalborg University, Aalborg, Denmark (E-mail: tbm@create.aau.dk).

problem. Extensive experimental results on the well-known OTB-2013, OTB-2015, TC-128, and VOT-2018 visual tracking datasets demonstrate that the proposed method effectively outperforms state-of-the-art methods in terms of precision and robustness of visual tracking.

**Keywords** Appearance modeling · discriminative correlation filters · deep convolutional neural networks · robust visual tracking

## 1 Introduction

Generic visual tracking aims to estimate the trajectory of an arbitrary visual target (given the initial state in the first frame) over time despite many challenging factors, including fast motion, background clutter, deformation, and occlusion. This constitutes a fundamental problem in many computer vision applications (e.g., self-driving cars, autonomous robots, human-computer interactions). Extracting a robust target representation is the critical component of state-of-the-art visual tracking methods to overcome these challenges. Hence, to robustly model target appearance, these methods utilize a wide range of handcrafted features (e.g., [41, 7, 87, 55, 92, 68] which exploit histogram of oriented gradients (HOG) [11], histogram of local intensities (HOI), and Color Names (CN) [81]), deep features from deep neural networks (e.g., [61, 63, 85, 2, 79, 40], or both (e.g., [14, 16, 64])).

Deep features are generally extracted from either FENs (i.e., pre-trained deep convolutional neural networks (CNNs)) or end-to-end networks (EENs), which directly evaluate target candidates [50]. However, most EEN-based visual trackers train or fine-tune FENs on visual tracking datasets. Numerous recent visual tracking methods exploit powerful generic target representations from FENs trained on large-scale object recognition datasets, such as the ImageNet [72]. By combining deep discriminative features from FENs with efficient online learning formulations, visual trackers based on discriminative correlation filters (DCF) have achieved significant performance in terms of accuracy, robustness, and speed. Despite the existence of modern CNN models, which are mainly based on the ResNet architectures [31], most DCF-based visual trackers still use VGG-M [5], VGG-16 [73], and VGG-19 [73] models, which have a simple stacked multi-layer topology with moderate representation capabilities. These models provide a limited power of representation compared to the state-of-the-art ResNet-based models such as ResNet [31], ResNeXt [90], squeeze-and-excitation networks (SE Nets) [37, 38], and DenseNet [39] models. Besides, most of the visual tracking methods exploit well-known CNN models which are trained for the visual recognition task. However, the transferability of deep representations from other related tasks (e.g., semantic segmentation or object detection) may help the DCF-based visual trackers to improve their robustness. Motivated by these two goals, the proposed paper surveys the performance of state-of-the-art ResNet-based FENs and exploits deep multitasking representations for visual tracking purposes.



The main contributions of the paper are as follows.

- 1) The effectiveness of twelve state-of-the-art ResNet-based FENs is evaluated: ResNet-50 [31], ResNet-101 [31], ResNet-152 [31], ResNeXt-50 [90], ResNeXt-101 [90], SE-ResNet-50 [37, 38], SE-ResNet-101 [37, 38], SE-ResNeXt-50 [37, 38], SE-ResNeXt-101 [37, 38], DenseNet-121 [39], DenseNet-169 [39], and DenseNet-201 [39]. These have been trained on the large-scale ImageNet dataset. To the best of our knowledge, this is the first paper to comprehensively survey the performance of ResNet-based FENs for visual tracking purposes.
- 2) The generalization of the best ResNet-based FEN is investigated by a different DCF-based visual tracking method.
- 3) A visual tracking method is proposed that fuses deep representations; these are extracted from the best trained ResNet-based network and the FCN-8s network [97], which have been trained on the PASCAL VOC [19] and MSCOCO [56] datasets for semantic image segmentation tasks.
- 4) The proposed method uses the extracted deep features from the FCN-8s network to semantically weight the target representation in each video frame.
- 5) The performance of the proposed method is extensively compared with state-of-the-art visual tracking methods on four well-known visual tracking datasets.

The rest of the paper is organized as follows. In Section 2, an overview of related work, including DCF-based visual trackers exploiting deep features from the FENs, is outlined. In Section 3 and Section 4, the survey of ResNet-based FENs and the proposed visual tracking method are presented, respectively. In Section 5, extensive experimental results on the large visual tracking datasets are given. Finally, the conclusion is summarized in Section 6.

## 2 Related Work

In this section, the exploitation of different FENs by state-of-the-art visual trackers is categorized and described. Although some visual tracking methods use recent FENs (e.g., R-CNN [25] in [35]), popular FENs in the most visual trackers are VGG-M, VGG-16, and VGG-19 models (see Table 1). These networks provide moderate accuracy and complexity for visual target modeling.

**VGG-M Model:** With the aid of deep features and spatial regularization weights, the spatially regularized discriminative correlation filters (DeepSRDCF) [14] learn more discriminative target models on larger image regions. The formulation of DeepSRDCF provides a larger set of negative training samples by penalizing the unwanted boundary effects resulting from the periodic assumption of standard DCF-based methods. Moreover, the method based on deep spatial-temporal regularized correlation filters (DeepSTRCF) [48] incorporates both spatial and temporal regularization parameters to provide a more robust appearance model. Then, it iteratively optimizes three closed-form solutions via the alternating direction method of multipliers (ADMM) algorithm [3]. To explore efficiency of Tikhonov regularization in tem-

**Table 1** Exploited FENS in some visual tracking methods.

Visual Tracking Method	Model	Pre-training Dataset	Name of Exploited Layer(s)
DeepSRDCF [14]	VGG-M	ImageNet	Conv1
C-COT [15]	VGG-M	ImageNet	Conv1, Conv5
ECO [16]	VGG-M	ImageNet	Conv1, Conv5
DeepSTRCF [48]	VGG-M	ImageNet	Conv3
WAEF [71]	VGG-M	ImageNet	Conv1, Conv5
WECO [32]	VGG-M	ImageNet	Conv1, Conv5
VDSR-SRT [57]	VGG-M	ImageNet	Conv1, Conv5
RPCF [77]	VGG-M	ImageNet	Conv1, Conv5
DeepTACF [47]	VGG-M	ImageNet	Conv1
ASRCF [10]	VGG-M, VGG-16	ImageNet	Norm1, Conv4-3
DRT [76]	VGG-M, VGG-16	ImageNet	Conv1, Conv4-3
ETDL [91]	VGG-16	ImageNet	Conv1-2
FCNT [82]	VGG-16	ImageNet	Conv4-3, Conv5-3
Tracker [26]	VGG-16	ImageNet	Conv2-2, Conv5-3
DNT [8]	VGG-16	ImageNet	Conv4-3, Conv5-3
CREST [74]	VGG-16	ImageNet	Conv4-3
CPT [6]	VGG-16	ImageNet	Conv5-1, Conv5-3
DeepFWDCE [18]	VGG-16	ImageNet	Conv4-3
DTO [86]	VGG-16, SSD	ImageNet	Conv3-3, Conv4-3, Conv5-3
DeepHPFT [51]	VGG-16, VGG-19, and GoogLeNet	ImageNet	Conv5-3, Conv5-4, and icp6-out
HCFT [61]	VGG-19	ImageNet	Conv3-4, Conv4-4, Conv5-4
HCFTs [63]	VGG-19	ImageNet	Conv3-4, Conv4-4, Conv5-4
LCTdeep [64]	VGG-19	ImageNet	Conv5-4
Tracker [62]	VGG-19	ImageNet	Conv3-4, Conv4-4, Conv5-4
HDT [70]	VGG-19	ImageNet	Conv4-2, Conv4-3, Conv4-4, Conv5-2, Conv5-3, Conv5-4
IBCCF [49]	VGG-19	ImageNet	Conv3-4, Conv4-4, Conv5-4
DCPF [65]	VGG-19	ImageNet	Conv3-4, Conv4-4, Conv5-4
MCPF [95]	VGG-19	ImageNet	Conv3-4, Conv4-4, Conv5-4
DeepLMCF [83]	VGG-19	ImageNet	Conv3-4, Conv4-4, Conv5-4
STSGS [94]	VGG-19	ImageNet	Conv3-4, Conv4-4, Conv5-4
MCCT [84]	VGG-19	ImageNet	Conv4-4, Conv5-4
DCPF2 [66]	VGG-19	ImageNet	Conv3-4, Conv4-4, Conv5-4
ORHF [58]	VGG-19	ImageNet	Conv3-4, Conv4-4, Conv5-4
IMM-DFT [78]	VGG-19	ImageNet	Conv3-4, Conv4-4, Conv5-4
MMLT [45]	VGGNet, Fully-convolutional Siamese network	ImageNet, ILSVRC-VID	Conv5
CNN-SVM [35]	R-CNN	ImageNet	First fully-connected layer
TADT [52]	Siamese matching network	ImageNet	Conv4-1, Conv4-3

poral domain, the weighted aggregation with enhancement filter (WAEF) [71] organizes deep features according to the average information content and suppresses uncorrelated frames for visual tracking. The continuous convolution operator tracker (C-COT) [15] fuses multi-resolution deep feature maps to learn discriminative continuous-domain convolution operators. It also exploits an implicit interpolation model to enable accurate sub-pixel localization of the visual target. To reduce the computational complexity and the number of training samples and to improve the update strategy of the C-COT, the ECO tracker [16] uses a factorized convolution operator, a compact generative model of training sample distribution, and a conservative model update strategy. The weighted ECO [32] utilizes weighted sum operation and normalization of deep features to take advantages of multi-resolution deep features. By employing the ECO framework, the VDSR-SRT method [57] uses a super-resolution reconstruction algorithm to robustly track targets in low-resolution images. To compress model size and improve the robustness against deformations, the region of interest (ROI) pooled correlation filters (RPCF) [77] accurately localizes a target while it uses deep feature maps with smaller sizes. The target-aware correlation filters (TACF) method [47] guides the learned filters for visual tracking by focusing on a target and preventing from background information.

**VGG-16 Model:** By employing the DeepSRDCF tracker, Gladh et al. [26] have investigated the fusion of handcrafted appearance features (e.g., HOG and CN) with deep RGB and motion features in the DCF-based visual tracking framework. This method demonstrates the effectiveness of deep feature maps for visual tracking purposes. Besides, the CREST method [74] reformulates correlation filters to extract more beneficial deep features by integrating

the feature extraction and learning process of DCFs. The dual network-based tracker (DNT) [8] designs a dual structure for embedding the boundary and shape information into the feature maps to better utilize deep hierarchical features. The deep fully convolutional networks tracker (FCNT) [82] uses two complementary feature maps and a feature-map selection method to design an effective FEN-based visual tracker. It exploits two different convolutional layers for category detection and distraction separation. Moreover, the feature-map selection method helps the FCNT to reduce computation redundancy and discard irrelevant feature maps. The method based on enhanced tracking and detection learning (ETDL) [91] employs different color bases for each color frame, adaptive multi-scale DCF with deep features, and a detection module to re-detect the target in failure cases. To prevent unexpected background information and distractors, the adaptive feature weighted DCF (FWDCF) [18] calculates target likelihood to provide spatial-temporal weights for deep features. The channel pruning method (CPT) [6] utilizes average feature energy ratio method to exploit low-dimensional deep features, adaptively.

**VGG-19 Model:** Similar to the FCNT, the hierarchical correlation feature-based tracker (HCFT) [61] exploits multiple levels of deep feature maps to handle considerable appearance variation and simultaneously to provide precise localization. Furthermore, the modified HCFT (namely HCFTs or HCFT\*) [63] not only learns linear correlation filters on multi-level deep feature maps; it also employs two types of region proposals and a discriminative classifier to provide a long-term memory of target appearance. Furthermore, the IMM-DFT method [78] exploits adaptive hierarchical features to interactively model a target due to the insufficiency of linear combination of deep features. To incorporate motion information with spatial deep features, the STSGS method [94] solves a compositional energy optimization to effectively localize an interested target. Also, the MCCT method [84] learns different target models by multiple DCFs that adopt different features, and then it makes the decision based on the reliable result. The hedged deep tracker (HDT) [70] ensembles weak CNN-based visual trackers via an online decision-theoretical Hedge algorithm, aiming to exploit the full advantage of hierarchical deep feature maps. The ORHF method [58] selects useful deep features according to the estimated confidence scores to reduce the computational complexity problem. To enhance the discrimination power of target among its background, the DCPF [65] and DeepLMCF [83] exploit deep features into the particle filter and structured SVM based tracking, respectively. To handle significant appearance change and scale variation, the deep long-term correlation tracking (LCTdeep) [64] uses multiple adaptive DCF, deep feature pyramid reconfiguration, aggressive and conservative learning rates as short- and long-term memories combined with an incrementally learned detector to recover tracking failures. Furthermore, Ma et al. [62] exploited different deep feature maps and a conservative update scheme to learn the mapping as a spatial correlation and maintain the long-memory of target appearance. At last, the DCPF2 [66] and IBCCF [49] methods aims to handle the aspect

ratio of a target for deep DCF-based trackers.

**Custom FENs:** In addition to the mentioned popular FENs, some methods exploit either a combination of FENs (e.g., ASRCF [10], DRT [76], DTO [86], MMLT [45], and DeepHPFT [51]) or other types of FENs (e.g., CNN-SVM [35], and TADT [52]) for visual tracking. The ASRCF [10] and DRT [76] methods modify the DCF formulation to achieve more reliable results. To alleviate inaccurate estimations and scale drift, the DeepHPFT [51] ensembles various handcrafted and deep features into the particle filter framework. Although the DTO method [86] employs a trained object detection model (i.e., single shot detector (SSD) [59]) to evaluate the impact of object category estimation for visual tracking, its ability is limited to some particular categories. The MMLT method [45] employs the Siamese and VGG networks to not only robustly track a target but also provide a long-term memory of appearance variation. On the other hand, the CNN-SVM [35] and TADT [52] methods use different models to improve the effectiveness of FENs for visual tracking.

According to the Table 1, the recent visual tracking methods do not only use the state-of-the-art FENs as their feature extractor; they also exploit different feature maps in the DCF-based framework. In the next section, the popular ResNet-based FENs will be surveyed. Then, the proposed visual tracking method, with the aid of fused deep representation, will be described.

### 3 Survey of ResNet-based FENs

This section has two main aims. First, twelve state-of-the-art ResNet-based FENs in the DCF-based framework are comprehensively evaluated for visual tracking purposes, and the best ResNet-based FEN and its feature maps are selected to be exploited in the DCF-based visual tracking framework. Second, the generalization of the best FEN is investigated on another DCF-based visual tracking method.

#### 3.1 Performance Evaluation

The ECO framework [16] is employed as the baseline visual tracker. Aiming for accurate and fair evaluations, this tracker was modified to extract deep features from FENs with the directed acyclic graph topology and exploit all deep feature maps without any down sampling or dimension reduction. Table 2 shows the exploited ResNet-based FENs and their feature maps, which are used in the proposed method. All these FENs are evaluated by well-known precision and success metrics [88, 89] on the OTB-2013 dataset [88], which includes more than 29,000 video frames. The characteristics of visual tracking datasets which are used in this work are shown in Table 3. The OTB-2013, OTB-2015, and TC-128 visual tracking datasets [88, 89, 54] have common challenging attributes, including illumination variation (IV), out-of-plane rotation (OPR), scale variation (SV), occlusion (OCC), deformation (DEF), motion blur (MB), fast motion (FM), in-plane rotation (IPR), out-of-view (OV), background clutter (BC), and low resolution (LR), which may occur for different classes of visual targets in real-world scenarios. Also, the

**Table 2** Exploited state-of-the-art ResNet-based FENs and their feature maps [Level of feature map resolution is denoted as L].

Name of FENs	Name of Output Layers (Number of Feature Maps)			
	L1	L2	L3	L4
ResNet-50	Conv1-relu (64)	Res2c-relu (256)	Res3d-relu (512)	Res4f-relu (1024)
ResNet-101	Conv1-relu (64)	Res2c-relu (256)	Res3b3-relu (512)	Res4b22-relu (1024)
ResNet-152	Conv1-relu (64)	Res2c-relu (256)	Res3b7-relu (512)	Res4b35-relu (1024)
ResNeXt-50	Features-2 (64)	Features-4-2-id-relu (256)	Features-5-3-id-relu (512)	Features-6-5-id-relu (1024)
ResNeXt-101	Features-2 (64)	Features-4-2-id-relu (256)	Features-5-3-id-relu (512)	Features-6-22-id-relu (1024)
SE-ResNet-50	Conv1-relu-7x7-s2 (64)	Conv2-3-relu (256)	Conv3-4-relu (512)	Conv4-6-relu (1024)
SE-ResNet-101	Conv1-relu-7x7-s2 (64)	Conv2-3-relu (256)	Conv3-4-relu (512)	Conv4-23-relu (1024)
SE-ResNeXt-50	Conv1-relu-7x7-s2 (64)	Conv2-3-relu (256)	Conv3-4-relu (512)	Conv4-6-relu (1024)
SE-ResNeXt-101	Conv1-relu-7x7-s2 (64)	Conv2-3-relu (256)	Conv3-4-relu (512)	Conv4-23-relu (1024)
DenseNet-121	Features-0-relu0 (64)	Features-0-transition1-relu (256)	Features-0-transition2-relu (512)	Features-0-transition3-relu (1024)
DenseNet-169	Features-0-relu0 (64)	Features-0-transition1-relu (256)	Features-0-transition2-relu (512)	Features-0-transition3-relu (1024)
DenseNet-201	Features-0-relu0 (64)	Features-0-transition1-relu (256)	Features-0-transition2-relu (512)	Features-0-transition3-relu (1792)

**Table 3** Exploited visual tracking datasets in this work.

Visual Tracking Dataset	Number of Videos	Number of Frames	Number of Videos Per Attribute											
			TV	OPR	SV	OCC	DEF	MB	FM	IPR	OV	BC	LR	
OTB-2013 [88]	51	29491	25	39	28	29	19	12	17	31	6	21	4	
OTB-2015 [89]	100	59040	38	63	65	49	44	31	40	53	14	31	10	
TC-128 [54]	129	55346	37	73	66	64	38	35	53	59	16	46	21	
VOT-2018 [43]	70	25504	Frame-based attributes											

VOT-2018 includes six main challenging attributes of camera motion, illumination change, motion change, occlusion, and size change which are annotated per each video frame. To measure visual tracking performance, the OTB-2013, OTB-2015, and TC-128 toolkits use precision and success plots to rank the methods according to the area under curve (AUC) metric while the VOT-2018 utilizes the accuracy-robustness (AR) plot to rank the visual trackers based on the TraX protocol [4]. The precision metric is defined as the percentage of frames where the average Euclidean distance between the estimated and ground-truth locations is smaller than a given threshold (20 pixels in this work). Moreover, the overlap success metric is the percentage of frames where the average overlap score between the estimated and the ground-truth bounding boxes is more than a particular threshold (50% overlap in this work). For the VOT-2018, the accuracy measure the overlap of the estimated bounding boxes with the ground-truth ones while the number of failures is considered as the robustness metric. Although the OTB-2013, OTB-2015, and TC-128 toolkits assess the visual trackers according to one-pass evaluation, the VOT-2014 detects the tracking failures of each method and restart the evaluations for each method after five frames of failure on every video sequence based on the TraX protocol.

To survey the performance of twelve ResNet-based FENs for visual tracking, the results of the comprehensive precision and success evaluations on the OTB-2013 dataset are shown in Fig. 1. On this basis, the L3 feature maps of these FENs have provided the best representation of targets, which is favorable for visual tracking. The deep features in the third level of ResNet-based FENs provide an appropriate balance between semantic information and spatial resolution and yet are invariant to significant appearance variations. The performance comparison of different ResNet-based FENs in terms of precision and success metrics is shown in Table 4. According to the results, the DenseNet-201 model and the feature maps extracted from the L3 layers have achieved the best performance for visual tracking purposes. Fur-

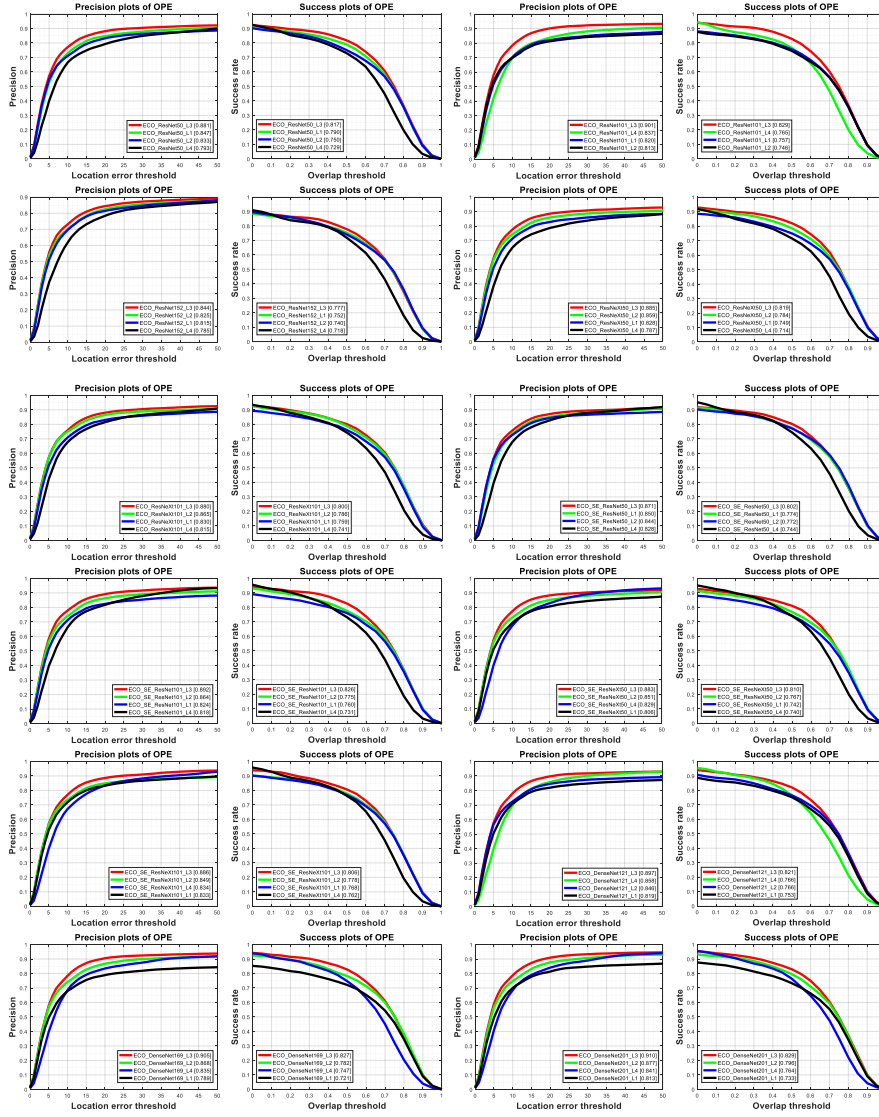


Fig. 1 Precision and success evaluation results of state-of-the-art ResNet-based FENs on the OTB-2013 dataset.

thermore, the comprehensive results show that the DenseNet-201 network provides superior deep representations of target appearance in challenging scenarios compared to the other ResNet-based FENs. This network features reuse property, which concatenates the feature maps learned by different layers to strengthen feature propagation and improve the efficiency. While this network is selected as the FEN in the proposed method (see Sec. 4), the ex-

**Table 4** Performance comparison of the L3 representations of ResNet-based models [first, second, and third FENs are shown in color].

Name of FENs	Precision	Success	Average
ResNet-50	0.881	0.817	0.849
ResNet-101	0.901	0.829	0.865
ResNet-152	0.844	0.777	0.810
ResNeXt-50	0.885	0.819	0.852
ResNeXt-101	0.880	0.800	0.840
SE-ResNet-50	0.871	0.802	0.836
SE-ResNet-101	0.892	0.826	0.859
SE-ResNeXt-50	0.883	0.810	0.846
SE-ResNeXt-101	0.886	0.806	0.846
DenseNet-121	0.897	0.821	0.859
DenseNet-169	0.905	0.827	0.866
DenseNet-201	0.910	0.829	0.869

ploitation of its best feature maps will be investigated in the next subsection.

### 3.2 Selection and Generalization Evaluation of Best Feature Maps

Two aims are investigated in the following section; exploiting the best feature maps of the DenseNet-201 model and performing generalization evaluation of this network in other DCF-based visual trackers.

In the first step, all single and fused deep feature maps extracted from the DenseNet-201 model were extensively evaluated on the OTB-2013 dataset (see Fig. 2(a)). Based on these results, the feature maps extracted from the L3 layer appeared to have achieved the best visual tracking performance in terms of average precision and success metrics. Although the original ECO tracker exploits fused deep features extracted from the first and third convolutional layers of the VGG-M network, the fusion of different levels of feature maps from the DenseNet network did not lead to better visual tracking performance. These results may help the visual trackers to prevent redundant feature maps and considerably reduce the computational complexity.

In the second step, the generalization of exploiting the third convolutional layer of the DenseNet-201 model was evaluated on the DeepSTRCF tracker [48]. To ensure a fair comparison, the visual tracking performance of both the original and the proposed versions of DeepSTRCF were evaluated with the aid of only deep features (without fusing with handcrafted features) extracted from VGG-M and DenseNet201 models on the OTB-2013 dataset. As shown in Fig. 2(b), the exploitation of the DenseNet-201 network does not only significantly improve the accuracy and the robustness of visual tracking in challenging scenarios; it can also generalize well into other DCF-based visual trackers.

## 4 Proposed Visual Tracking Method

The proposed tracking method exploits fused deep features of multi-task FENs and introduces semantic weighting to provide a more robust target representation for the learning process. To effectively track generic objects in challenging scenarios, the proposed method aims to exploit fused deep representation (extracted from FENs) and semantic weighting of target regions. The algorithm 1 is shown a brief outline of the proposed method.

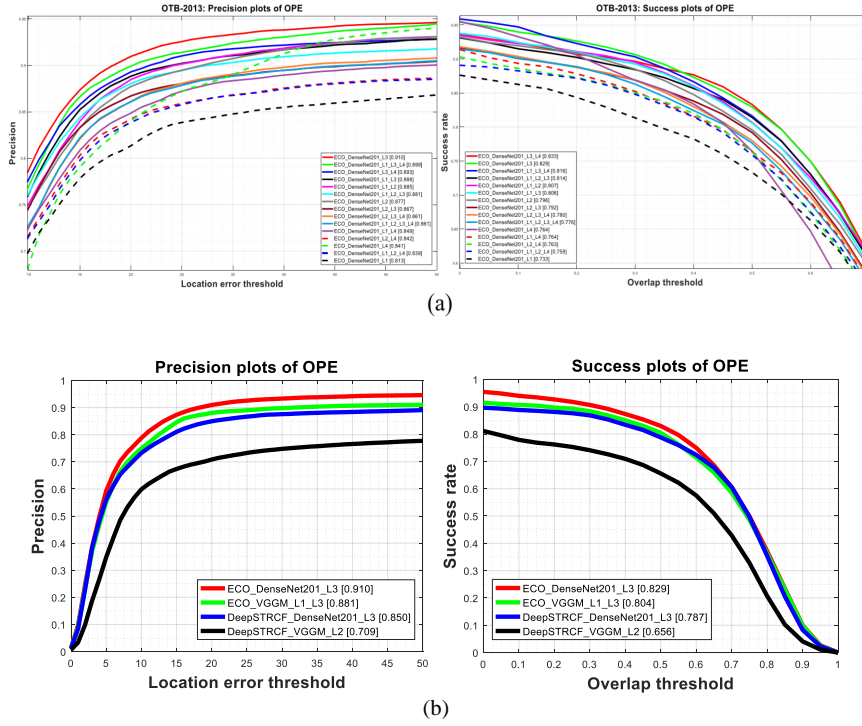


Fig. 2 Evaluation of exploiting fused feature maps of the pre-trained DenseNet-201 network and generalization of the best of them into the DeepSTRCF tracker.

#### Algorithm 1. Brief Outline of Proposed Method

##### – Pre-requisitions:

- Survey of twelve ResNet-based models which have been pre-trained on ImageNet dataset
- Selection of the best feature maps of the DenseNet-201 model and their generalization
- Selection of the DenseNet-201 model as the best FEN
- Selection of the FCN-8s model which has been pre-trained on PASCAL VOC and MSCOCO datasets
- Selection of the pixel-wise prediction maps as the semantic features

##### – Evaluation of proposed visual tracking method on the OTB-2013, OTB-2015, and TC-128 datasets (without any pre-training of translation and scale filters):

Initialization of the proposed method with a given bounding box in the first frame:

1. Extract deep features from DenseNet-201 and FCN-8s models
2. Fuse deep multi-tasking features with semantic windowing using Eq. (2) to Eq. (7)
3. Learn the continuous translation filters using Eq. (9) to robustly model the target
4. Extract HOG features
5. Learn continuous scale filters using a multi-scale search strategy
6. Update the translation and scale filters using Eq. (11) (for frame  $t > 1$ )
7. Go to the next frame
8. Select the search region according to the previous target location
9. Extract deep features from DenseNet-201 and FCN-8s models
10. Fuse deep multi-tasking features with semantic windowing using Eq. (2) to Eq. (7)
11. Calculate the confidence map Eq. (9)
12. Estimate the target location by finding the maximum of Eq. (10)
13. Extract multi-scale HOG features centered at the estimated target location
14. Estimate the target scale using the multi-scale search strategy
15. Select the target region conforming with the estimated location and scale
16. Return to Step 1



Using the ECO framework, the proposed method minimizes the following objective function in the Fourier domain [16]

$$\min_H \sum_{i=1}^m \alpha_i \left\| \sum_{j=1}^n H^j F_i^j K_j - Y_i \right\|_{l^2}^2 + \sum_{j=1}^n \|P * H^j\|_{l^2}^2 \quad (1)$$

in which  $H$ ,  $F$ ,  $Y$ , and  $K$  denote the multi-channel convolution filters, weighted fused deep feature maps from FENs, desirable response map, and interpolation function (to pose the learning problem in the continuous domain), respectively. The penalty function (to learn filters on target region), number of training samples, number of convolution filters, and weights of training samples are indicated by the  $P$ ,  $m$ ,  $n$ , and  $\alpha_i$  variables, respectively. Also, the  $*$ , and  $l^2$  are referred to the circular convolution operation and L2 norm, respectively. Capital letters represent the discrete Fourier transform (DFT) of variables.

In addition to utilizing the feature maps of the DenseNet-201 network, the proposed method extracts deep feature maps from the FCN-8s network (i.e., 21 feature maps from the “*score\_final*” layer) to learn a more discriminative target representation by fusing the feature maps as

$$\bar{X} = [X_{DenseNet}^1, \dots, X_{DenseNet}^{S1}, X_{FCN}^1, \dots, X_{FCN}^{S2}] \quad (2)$$

where  $\bar{X}$ ,  $X$ ,  $S1$ , and  $S2$  are the fused deep representation, deep feature maps, number of feature maps extracted from the DenseNet-201 network, and number of feature maps extracted from the FCN-8s network, respectively.

Although the DCF-based visual trackers use a cosine window to weight the target region and its surroundings, this weighting scheme may contaminate the target model and lead to drift problems. Some recent visual trackers, e.g., [44], utilize different windowing methods, such as computation of target likelihood maps from raw pixels. However, the present work is the first method to use deep feature maps to semantically weight target appearances. The proposed method exploits the feature maps of the FCN-8s model to semantically weight the target. Given the location of the target at  $(e, f)$  in the first (or previous) frame, the proposed method defines a semantic mask  $W_{mask}$  as

$$L_{FCN}(i, j) = \begin{cases} Z & \text{if } \max X_{FCN}^Z(i, j) > \max_{C \neq Z} X_{FCN}^C(i, j) \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

$$W_{mask}(i, j) = \begin{cases} 1 & \text{if } L_{FCN}(i, j) = L_{FCN}(e, f) \\ 0 & \text{otherwise} \end{cases} \quad (4)$$

in which  $L_{FCN}$ ,  $Z$ , and  $(i, j)$  are the label map, PASCAL classes  $\{C, Z\} \in (1, \dots, 21)$ , and spatial location, respectively. Finally, the fused feature maps are weighted by

$$F = \bar{X} \odot (W_{mask} \cdot W_{cos}) \quad (5)$$

where  $W_{cos}$  and  $\odot$  are the conventional cosine (or sine) window and the channel-wise product, respectively. The traditional cosine window weights a deep representation  $X$  with  $d1$  and  $d2$  dimensions as

$$X_{ij}^w = (X_{ij} - 0.5) \sin(\pi i/d1) \sin(\pi j/d2), \forall i = 0, \dots, d1 - 1, j = 0, \dots, d2 - 1 \quad (6)$$

The matrix form of semantically weighted fused deep representations is defined as

$$B \triangleq FK = [\bar{X} \odot (W_{mask} \cdot W_{cos})] \cdot K \quad (7)$$

such that the weighted fused feature maps are posed to the continuous domain by an interpolation function. By defining the diagonal matrix of training samples weights as  $\Gamma$ , the regularized least square problem (1) is minimized by

$$\nabla_H \|B\Gamma H - Y\|_2^2 + \|PH\|_2^2 = 0 \quad (8)$$

which has the following closed-form solution

$$[B^T \Gamma B + P^T P] H = B^T \Gamma Y \quad (9)$$

The proposed method learns multi-channel convolution filters by employing the iterative conjugate gradient (CG) method. Because the iterative CG method does not need to form the multiplications of  $B^T$  and  $B$ , it can desirably minimize the objective function in limited iterations. To track the target in subsequent frames, the multitasking deep feature maps ( $T$ ) are extracted from the search region in the current frame. Next, the confidence map is calculated by

$$CM = F^{-1} \left\{ \sum_{j=1}^n T^j \cdot (H^j F^j K_j) \right\} \quad (10)$$

that its maximum determines the target location in the current frame. Then, a new set of continuous convolution filters are trained on the current target sample and the prior learned filters are updated by the current learned filters as

$$H = (1 - \gamma) H_{prior} + \gamma H_{current} \quad (11)$$

in which  $\gamma$  is learning rate. Therefore, the proposed method alternatively uses online training and tracking processes of the generic visual target for all video frames.

## 5 Experimental Results

All implementation details related to the exploited models (including the ResNet, ResNeXt, SE-ResNet, SE-ResNeXt, DenseNet, and FCN-8s models) are the same as in their original papers [31, 90, 37, 38, 39, 97]. The ResNet-based models have been trained on the ImageNet dataset [72] while the FCN-8s model has been trained on the PASCAL VOC and MSCOCO data-sets [19, 56]. Similar to the exploitation of FENs in state-of-the-art visual tracking methods [14, 26, 48, 15, 16, 91, 82, 61, 63, 8, 70, 64, 62], these CNNs are

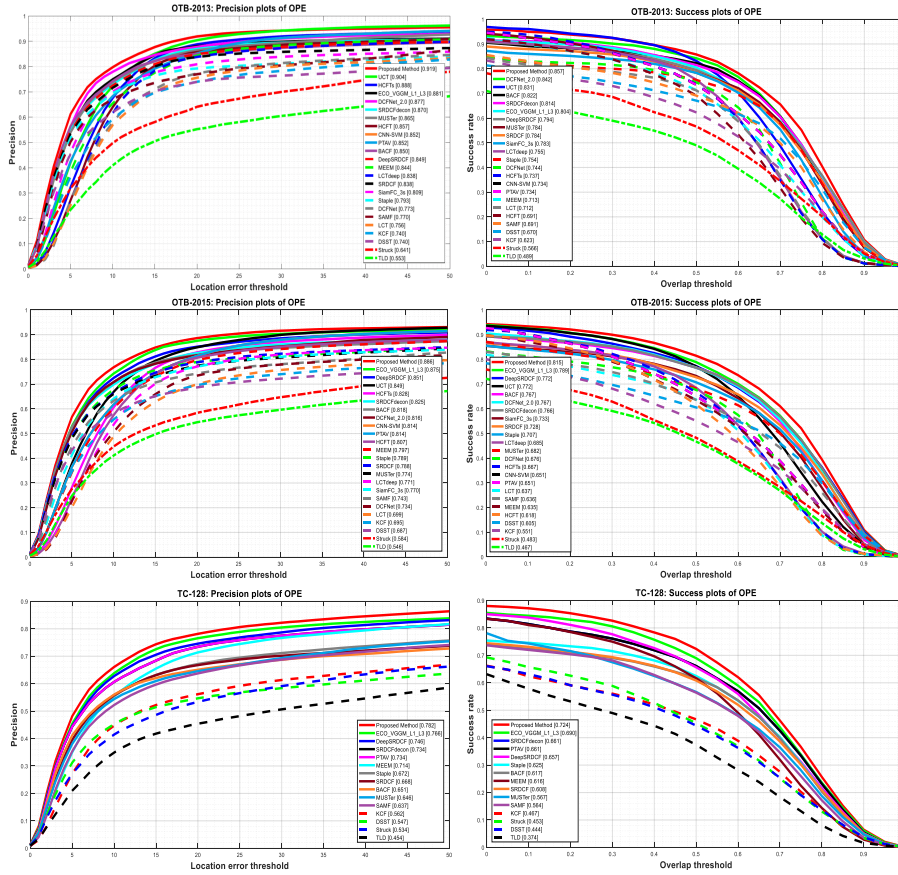
not fine-tuned or trained during their performance evaluations on the OTB-2013, OTB-2015, TC-128, and VOT-2018 visual tracking datasets. In fact, the FENs transfers the power of deep features to the visual tracking methods. According to the visual tracking benchmarks, all visual trackers are initialized with only a given bounding box in the first frame of each video sequence. Then, these visual trackers start to learn and track a target, simultaneously. To investigate the effectiveness of the proposed method and fair evaluations, the experimental configuration of baseline ECO tracker is used. The search region, learning rate, number of CG iterations, and number of training samples are set to the  $5^2$  times of the target box, 0.009, 5, and 50, respectively. The desirable response map is considered as a 2D Gaussian function which its standard deviation is set to 0.083. Given a target with size  $q_1 \times q_2$  and spatial location  $(i, j)$ , the penalty function is defined as

$$p(i, j) = 0.1 + 3 \left( \frac{i}{q_1} \right)^2 + 3 \left( \frac{j}{q_2} \right)^2 \quad (12)$$

which is a quadratic function with a minimum at its center. Moreover, the number of deep feature maps extracted from the DenseNet and FCN-8s models are  $S_1 = 512$  and  $S_2 = 21$ . To estimate the target scale, the proposed method utilizes the multi-scale search strategy [17] with 17 scales which have relative scale factor 1.02. Note that this multi-scale search strategy exploits the HOG features to accelerate the visual tracking methods (such as in the DSST, C-COT, and ECO). Like the baseline ECO tracker, the proposed method applies the same configuration to all videos in different visual tracking datasets. However, this work exploits only deep features to model target appearance and highlight the effectiveness of the proposed method. This section includes a comparison of the proposed method with state-of-the-art visual tracking methods and the ablation study, which evaluates the effectiveness of proposed components on visual tracking performance.

### 5.1 Performance Comparison

To evaluate the performance of the proposed method, it is extensively compared quantitatively with deep and handcrafted-based state-of-the-art visual trackers (for which benchmark results on different datasets have been publicly available) in the one-pass evaluation (OPE) on three large visual tracking datasets: OTB-2013 [88], OTB-2015 [89], TC-128 [54], and VOT-2018 [43] datasets. The state-of-the-art deep visual trackers include deep feature-based visual trackers (namely, ECO-CNN [16], DeepSRDCF [14], UCT [98], DCFNet [85], DCFNet-2.0 [85], LCTdeep [64], HCFT [61], HCFTs [63], SiamFC-3s [2], CNN-SVM [35], PTAV [21, 20]), CCOT [15], DSiam [28], CFNet [80], DeepC-SRDCF [60], MCPF [95], TRACA [9], DeepSTRCF [48], SiamRPN [46], SA-Siam [30], LSART [75], DRT [76], DAT [69], CFCF [27], CRPN [22], GCT [24], SiamDW-SiamFC [96], and SiamDW-SiamRPN [96]; Also the compared handcrafted feature-based visual trackers are the Struck [29], BACF [23],



**Fig. 3** Overall precision and success evaluations on the OTB-2013, OTB-2015, and TC-128 datasets.

DSST [17], MUSTer [36], SRDCF [13], SRDCFdecon [12], Staple [1], LCT [64], KCF [33, 34], SAMF [53], MEEM [93], and TLD [42]). These methods are included the visual trackers that exploit handcrafted features, deep features (by FENs or EENs), or both. The proposed method is implemented on an Intel I7-6800K 3.40 GHz CPU with 64 GB RAM with an NVIDIA GeForce GTX 1080 GPU. The overall and attribute-based performance comparisons of the proposed method with different visual trackers on the four visual tracking datasets in terms of various evaluation metrics are shown in Fig. 3 to Fig. 5.

Based on the results of the attribute-based comparison on the OTB-2015 dataset, the proposed method has shown to improve the average success rate of the ECO tracker by up to 5.1%, 1.8%, 1.3%, 1.1%, 4.3%, 2.4%, 2.8%, and 6.5% for IV, OPR, SV, OCC, DEF, FM, IPR, and BC attributes, respectively. However, the ECO tracker has achieved better visual tracking performance for MB, OV, and LR attributes. The proposed method outperforms other state-

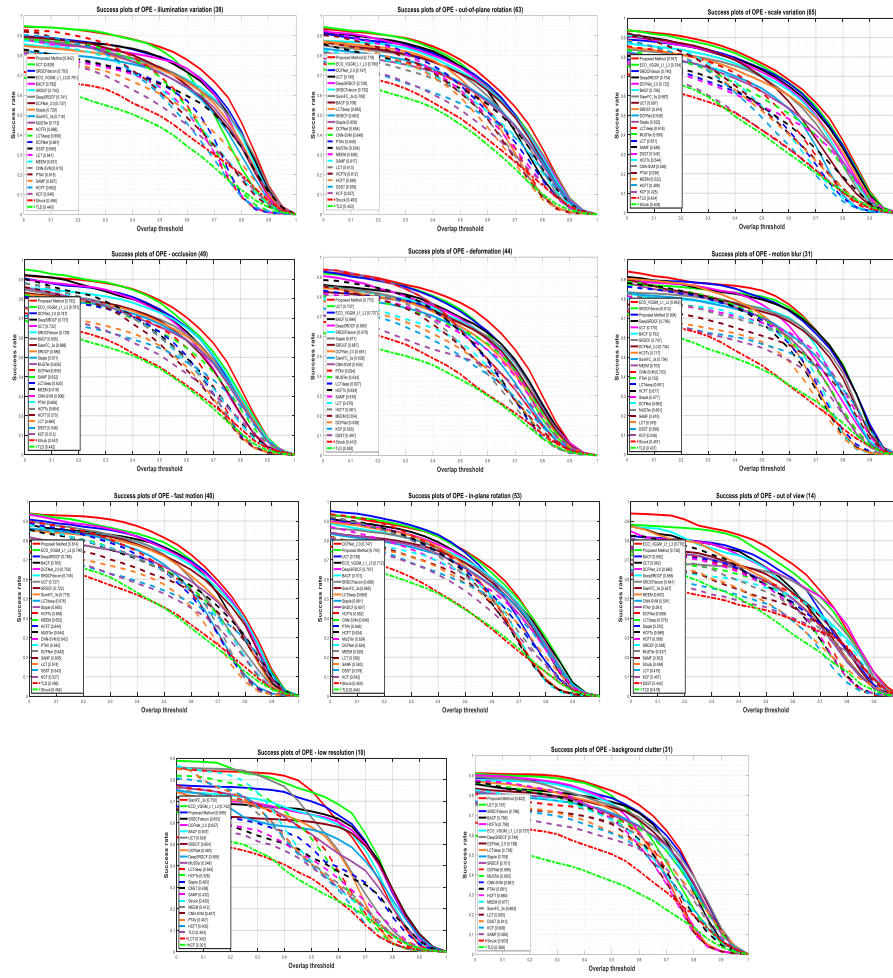


Fig. 4 Attribute-based evaluations of visual trackers in terms of average success rate on the OTB-2015 dataset.

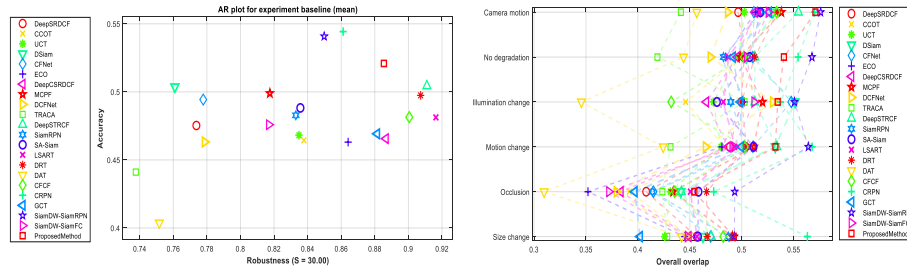


Fig. 5 Overall and attribute-based performance comparison of visual tracking methods on VOT-2018 dataset.

of-the-art visual trackers in most challenging attributes, and the results indicate that the proposed method has improved the average precision rate by up to 3.8%, 1.1%, and 1.6% and the average success rate by up to 5.3%, 2.6%, and 3.4% compared to the ECO-CNN tracker on the OTB-2013, OTB-2015, and TC-128 datasets, respectively. Thus, in addition to providing better visual tracking performance compared to the FEN-based visual trackers, the proposed method outperforms the EEN-based visual tracking methods. For example, the proposed method achieved up to 3.7%, 7%, 7.2%, 11.6%, and 15.2% on the OTB-2015 dataset in terms of average precision rate compared to the UCT, DCFNet-2.0, PTAV, SiamFC-3s, and DCFNet trackers, respectively. Furthermore, the success rate of the proposed method gained by up to 4.3%, 4.8%, 16.4%, 8.2%, and 13.9%, respectively, compared to the aforementioned EEN-based trackers. Because of the higher performance of deep visual trackers compared to the handcrafted-based ones, the evaluations on the VOT-2018 compare the proposed method just with the state-of-the-art deep visual tracking methods. The achieved results in Fig. 5 demonstrate the proficiency of the proposed method for visual tracking. According to these results, the proposed method not only considerably improve the performance of ECO-CNN framework but also has a competitive results compared to the EEN-based methods which have trained on large-scale datasets.

The qualitative comparison between the proposed method and the top-5 visual trackers (i.e., UCT, DCFNet-2.0, ECO-CNN, HCFTs, and DeepSRDCF) is shown in Fig. 5, which clearly demonstrates remarkable robustness of the proposed method in real-world scenarios. It is noteworthy that the proposed method provides an acceptable average speed (about five frames per second (FPS)) compared to the most FEN-based visual trackers, which run less than one FPS [67, 15]. To robustly model the visual targets, the proposed method uses two main components. First, the proposed method exploits the weighted fused deep representations in the learning process of convolution filters in the continuous domain. It extracts deep feature maps from two FENs (i.e., DenseNet-201 and FCN-8s) that are trained on different large-scale datasets (i.e., ImageNet, PASCAL VOC, and MSCOCO) for different tasks of object recognition and semantic segmentation. These rich representations provide complementary information for an unknown visual target. Hence, the fused deep representation leads to a more robust target model, which is more resist against challenging attributes, such as heavy illumination variation, occlusion, deformation, background clutter, and fast motion. Moreover, the employed FENs are efficient and do not have complications with integrating into the proposed DCF-based tracker. Second, the proposed method uses the extracted feature maps from the FCN-8s network to semantically enhance the weighting process of target representations. This component helps the proposed method to model the visual target more accurately and prevent contamination of target appearance with the background or visual distractors. Hence, the proposed method demonstrates more robustness against significant variations of a visual target, including deformation and background clutter.



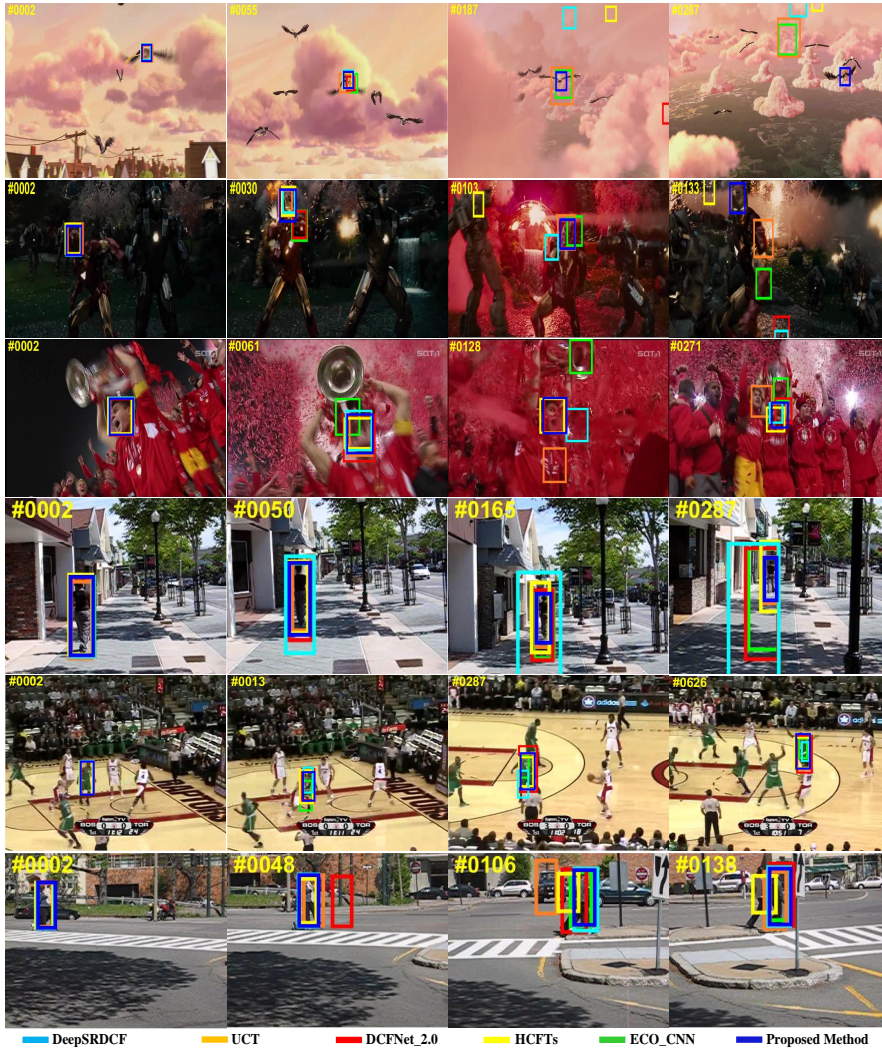


Fig. 6 Qualitative evaluation results of visual tracking methods on Bird1, Ironman, Soccer, Human9, Basketball, and Couple video sequences from top to bottom row, respectively.

## 5.2 Performance Comparison: Ablation Study

To understand the advantages and disadvantages of each proposed component, the visual tracking performance of two ablated trackers are investigated and compared to the proposed method. The ablated trackers are the proposed method to disable the fusion process (i.e., this method just exploits the DenseNet-201 representations and semantic windowing process) and the proposed method without semantic windowing process (i.e., this method only fuses deep multitasking representations from the DenseNet-201 and FCN-8s networks to ensure robust target appearance).

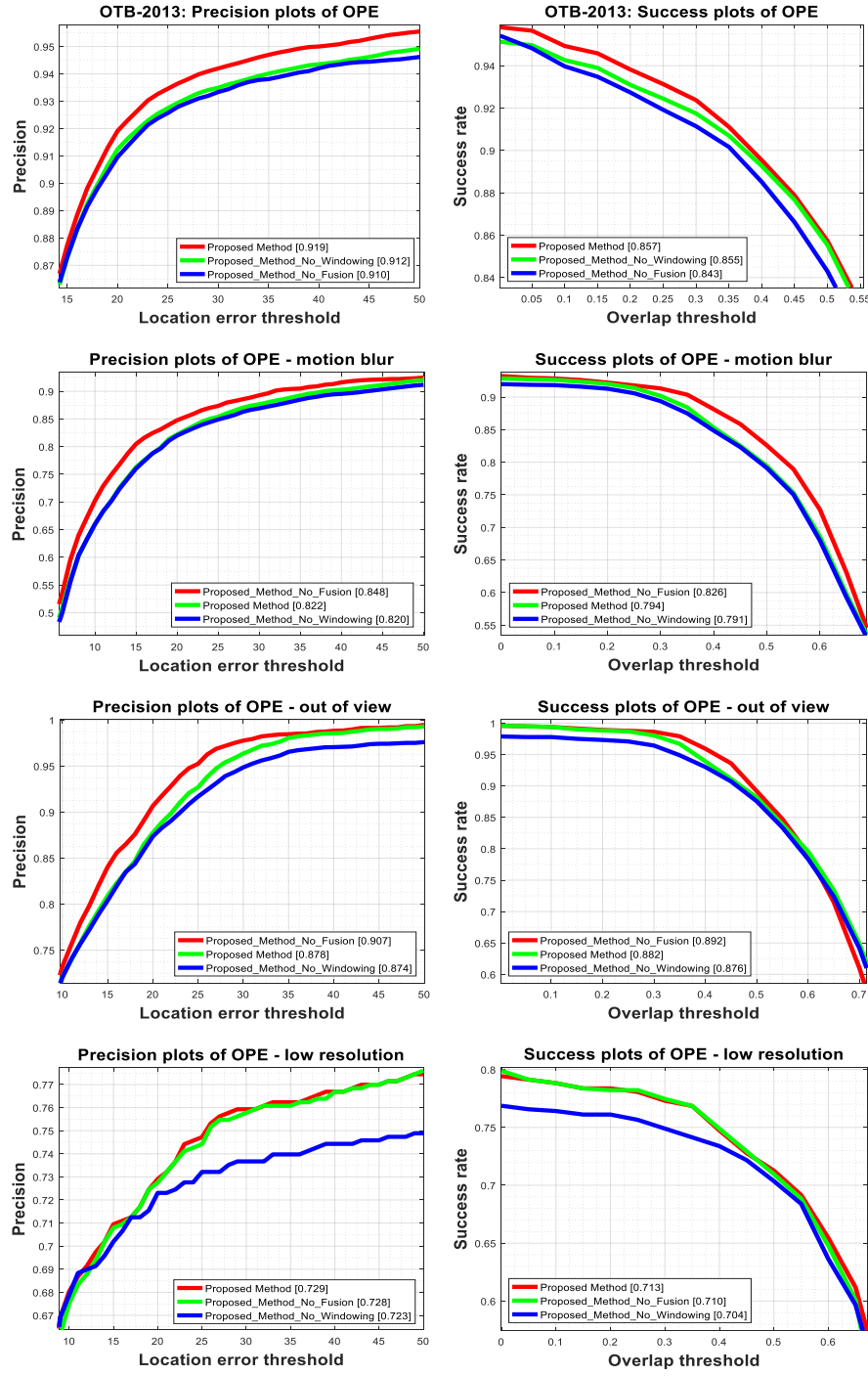


Fig. 7 Overall and attribute-based ablation study of proposed visual-tracking method on the OTB-2013 dataset.



The results of the ablation study on the OTB-2013 dataset are shown in Fig. 7. Based on the precision and success plots of the overall comparison, it can be concluded that most of the success of the proposed method is related to the effective fusion process of desirable representations for visual tracking. Moreover, it is clearly demonstrated that the proposed windowing process enhances the proposed method, particularly in terms of precision metrics. To identify the reason for the performance degradation faced for MB, OV, and LR attributes, an attribute-based evaluation was performed. The fusion process provides the main contribution to reducing visual tracking performance when these attributes occur, as seen in the results in Fig. 7. These results confirm that the semantic windowing process can effectively support visual trackers in learning more discriminative target appearance and in preventing drift problems.

## 6 Conclusion

The use of twelve state-of-the-art ResNet-based FENs for visual tracking purposes was evaluated. A visual tracking method was proposed; this method can fuse deep features extracted from the DenseNet-201 and FCN-8s networks, but it can also semantically weight target representations in the learning process of continuous convolution filters. Moreover, the best ResNet-based FEN in the DCF-based framework was determined (i.e., DenseNet-201), and the effectiveness of using the DenseNet-201 network on the tracking performance of another DCF-based tracker was explored. Finally, the comprehensive experimental results on the OTB-2013, OTB-2015, TC-128 and VOT-2018 datasets demonstrated that the proposed method outperforms state-of-the-art visual tracking methods in terms of various evaluation metrics for visual tracking.

**Acknowledgement:** This work was partly supported by a grant (No. 96013046) from Iran National Science Foundation (INSF).

**Compliance with Ethical Standards (Conflict of Interest):** All authors declare that they have no conflict of interest.

## References

1. Luca Bertinetto, Jack Valmadre, Stuart Golodetz, Ondrej Miksik, and Philip H.S. Torr. Staple: Complementary learners for real-time tracking. In *Proc. IEEE CVPR*, pages 1401–1409, 2016.
2. Luca Bertinetto, Jack Valmadre, João F. Henriques, Andrea Vedaldi, and Philip H.S. Torr. Fully-convolutional Siamese networks for object tracking. In *Proc. ECCV*, pages 850–865, 2016.
3. Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, and Jonathan Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in Machine Learning*, 3:1–122, 2010.

4. Luka Čehovin. TraX: The visual tracking exchange protocol and library. *Neurocomputing*, 260:5–8, 2017.
5. Ken Chatfield, Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Return of the devil in the details: Delving deep into convolutional nets. In *Proc. BMVC*, pages 1–11, 2014.
6. Manqiang Che, Runling Wang, Yan Lu, Yan Li, Hui Zhi, and Changzhen Xiong. Channel pruning for visual tracking. In *Proc. ECCVW*, pages 70–82, 2019.
7. Zhi Chen, Peizhong Liu, Yongzhao Du, Yanmin Luo, and Jing-Ming Guo. Robust visual tracking using self-adaptive strategy. *Multimed. Tools Appl.*, 2019.
8. Zhizhen Chi, Hongyang Li, Huchuan Lu, and Ming Hsuan Yang. Dual deep network for visual tracking. *IEEE Trans. Image Process.*, 26(4):2005–2015, 2017.
9. Jongwon Choi, Hyung Jin Chang, Tobias Fischer, Sangdoo Yun, Kyue-wang Lee, Jiyeoup Jeong, Yiannis Demiris, and Jin Young Choi. Context-aware deep feature compression for high-speed visual tracking. In *Proc. IEEE CVPR*, pages 479–488, 2018.
10. Kenan Dai, Dong Wang, Huchuan Lu, Chong Sun, and Jianhua Li. Visual tracking via adaptive spatially-regularized correlation filters. In *Proc. CVPR*, pages 4670–4679, 2019.
11. Navneet Dalal and Bill Triggs. Histograms of oriented gradients for human detection. In *Proc. IEEE CVPR*, pages 886–893, 2005.
12. M. Danelljan, G. Häger, F. S. Khan, and M. Felsberg. Adaptive decontamination of the training set: A unified formulation for discriminative visual tracking. In *Proc. IEEE CVPR*, pages 1430–1438, 2016.
13. Martin Danelljan, Gustav Hager, Fahad Shahbaz Khan, and Michael Felsberg. Learning spatially regularized correlation filters for visual tracking. In *Proc. IEEE ICCV*, pages 4310–4318, 2015.
14. Martin Danelljan, Gustav Hager, Fahad Shahbaz Khan, and Michael Felsberg. Convolutional features for correlation filter based visual tracking. In *Proc. IEEE ICCVW*, pages 621–629, 2016.
15. Martin Danelljan, Andreas Robinson, Fahad Shahbaz Khan, and Michael Felsberg. Beyond correlation filters: Learning continuous convolution operators for visual tracking. In *Proc. ECCV*, volume 9909 LNCS, pages 472–488, 2016.
16. Martin Danelljan, Goutam Bhat, Fahad Shahbaz Khan, and Michael Felsberg. ECO: Efficient convolution operators for tracking. In *Proc. IEEE CVPR*, pages 6931–6939, 2017.
17. Martin Danelljan, Gustav Hager, Fahad Shahbaz Khan, and Michael Felsberg. Discriminative Scale Space Tracking. *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(8):1561–1575, 2017.
18. Fei Du, Peng Liu, Wei Zhao, and Xianglong Tang. Spatial-temporal adaptive feature weighted correlation filter for visual tracking. *Signal Proc.: Image Comm.*, 67:58–70, 2018.

19. Mark Everingham, S. M. Ali Eslami, Luc Van Gool, Christopher K. I. Williams, John Winn, and Andrew Zisserman. The pascal visual object classes challenge: A retrospective. *IJCV*, 111(1):98–136, 2015.
20. Heng Fan and H.Ling. Parallel tracking and verifying. *IEEE Trans. Image Process.*, 28(8):4130–4144, 2019.
21. Heng Fan and Haibin Ling. Parallel tracking and verifying: A framework for real-time and high accuracy visual tracking. In *Proc. IEEE ICCV*, pages 5487–5495, 2017.
22. Heng Fan and Haibin Ling. Siamese cascaded region proposal networks for real-time visual tracking, 2018. URL <http://arxiv.org/abs/1812.06148>.
23. Hamed Kiani Galoogahi, Ashton Fagg, and Simon Lucey. Learning background-aware correlation filters for visual tracking. In *Proc. IEEE ICCV*, pages 1144–1152, 2017.
24. Junyu Gao, Tianzhu Zhang, and Changsheng Xu. Graph convolutional tracking. In *Proc. CVPR*, pages 4649–4659, 2019.
25. Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proc. IEEE CVPR*, pages 580–587, 2014.
26. Susanna Gladh, Martin Danelljan, Fahad Shahbaz Khan, and Michael Felsberg. Deep motion features for visual tracking. In *Proc. ICPR*, pages 1243–1248, 2016.
27. Erhan Gundogdu and A. Aydin Alatan. Good features to correlate for visual tracking. *IEEE Trans. Image Process.*, 27(5):2526–2540, 2018.
28. Qing Guo, Wei Feng, Ce Zhou, Rui Huang, Liang Wan, and Song Wang. Learning dynamic Siamese network for visual object tracking. In *Proc. IEEE ICCV*, pages 1781–1789, 2017.
29. Sam Hare, Stuart Golodetz, Amir Saffari, Vibhav Vineet, Ming Ming Cheng, Stephen L. Hicks, and Philip H.S. Torr. Struck: Structured output tracking with kernels. *IEEE Trans. Pattern Anal. Mach. Intell.*, 38(10):2096–2109, 2016.
30. Anfeng He, Chong Luo, Xinmei Tian, and Wenjun Zeng. A twofold Siamese network for real-time object tracking. In *Proc. IEEE CVPR*, pages 4834–4843, 2018.
31. Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proc. IEEE CVPR*, pages 770–778, 2016.
32. Zhiqun He, Yingruo Fan, Junfei Zhuang, Yuan Dong, and Hongliang Bai. Correlation filters with weighted convolution responses. In *Proc. ICCVW*, pages 1992–2000, 2018.
33. João F. Henriques, Rui Caseiro, Pedro Martins, and Jorge Batista. Exploiting the circulant structure of tracking-by-detection with kernels. In *Proc. ECCV*, pages 702–715, 2012.
34. Joao F. Henriques, Rui Caseiro, Pedro Martins, and Jorge Batista. High-speed tracking with kernelized correlation filters. *IEEE Trans. Pattern Anal. Mach. Intell.*, 37(3):583–596, 2015.

35. Seunghoon Hong, Tackgeun You, Suha Kwak, and Bohyung Han. On-line tracking by learning discriminative saliency map with convolutional neural network. In *Proc. ICML*, pages 597–606, 2015.
36. Zhibin Hong, Zhe Chen, Chaohui Wang, Xue Mei, Danil Prokhorov, and Dacheng Tao. MUlti-Store Tracker (MUSTer): A cognitive psychology inspired approach to object tracking. In *Proc. IEEE CVPR*, pages 749–758, 2015.
37. J. Hu, L. Shen, and G. Sun. Squeeze-and-excitation networks. In *Proc. IEEE CVPR*, pages 7132–7141, 2018.
38. J. Hu, L. Shen, S. Albanie, G. Sun, and E. Wu. Squeeze-and-excitation networks. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2019. doi: 10.1109/TPAMI.2019.2913372.
39. G. Huang, Z. Liu, L. v. d. Maaten, and K. Q. Weinberger. Densely connected convolutional networks. In *Proc. IEEE CVPR*, pages 2261–2269, 2017.
40. W. Huang, J. Gu, and X. et al. Ma. End-to-end multitask Siamese network with residual hierarchical attention for real-time object tracking. *Appl. Intell.*, 2020. URL <https://doi.org/10.1007/s10489-019-01605-2>.
41. Yang Huang, Zhiqiang Zhao, Bin Wu, Zhuolin Mei, Zongmin Cui, and Guangyong Gao. Visual object tracking with discriminative correlation filtering and hybrid color feature. *Multimed. Tools Appl.*, 2019.
42. Z. Kalal, K. Mikolajczyk, and J. Matas. Tracking-learning-detection. *IEEE Trans. Pattern Anal. Mach. Intell.*, 34(7):1409–1422, 2012.
43. Matej Kristan, Aleš Leonardis, Jiří Matas, Michael Felsberg, Roman Pflugfelder, and et al. The sixth visual object tracking vot2018 challenge results. In *Proc. ECCVW*, pages 3–53, 2019.
44. Yangliu Kuai, Gongjian Wen, and Dongdong Li. Learning adaptively windowed correlation filters for robust tracking. *J. VIS. COMMUN. IMAGE R.*, 51:104 – 111, 2018.
45. Hankyeol Lee, Seokeon Choi, and Changick Kim. A memory model based on the Siamese network for long-term tracking. In *Proc. ECCVW*, pages 100–115, 2019.
46. Bo Li, Junjie Yan, Wei Wu, Zheng Zhu, and Xiaolin Hu. High performance visual tracking with siamese region proposal network. In *Proc. IEEE CVPR*, pages 8971–8980, 2018.
47. Dongdong Li, Gongjian Wen, Yangliu Kuai, Jingjing Xiao, and Fatih Porikli. Learning target-aware correlation filters for visual tracking. *J. VIS. COMMUN. IMAGE R.*, 58:149–159, 2019.
48. Feng Li, Cheng Tian, Wangmeng Zuo, Lei Zhang, and Ming Hsuan Yang. Learning spatial-temporal regularized correlation filters for visual tracking. In *Proc. IEEE CVPR*, pages 4904–4913, 2018.
49. Feng Li, Yingjie Yao, Peihua Li, David Zhang, Wangmeng Zuo, and Ming Hsuan Yang. Integrating boundary and center correlation filters for visual tracking with aspect ratio variation. In *Proc. IEEE ICCVW*, pages 2001–2009, 2018.

50. Peixia Li, Dong Wang, Lijun Wang, and Huchuan Lu. Deep visual tracking: Review and experimental comparison. *Pattern Recognit.*, 76:323–338, 2018.
51. Shengjie Li, Shuai Zhao, Bo Cheng, Erhu Zhao, and Junliang Chen. Robust visual tracking via hierarchical particle filter and ensemble deep features. *IEEE Trans. Circuits Syst. Video Technol.*, 2018.
52. Xin Li, Chao Ma, Baoyuan Wu, Zhenyu He, and Ming-Hsuan Yang. Target-aware deep tracking, 2019. URL <http://arxiv.org/abs/1904.01772>.
53. Yang Li and Jianke Zhu. A scale adaptive kernel correlation filter tracker with feature integration. In *Proc. ECCVW*, pages 254–265, 2015.
54. Pengpeng Liang, Erik Blasch, and Haibin Ling. Encoding color information for visual tracking: Algorithms and benchmark. *IEEE Trans. Image Process.*, 24(12):5630–5644, 2015.
55. Yun Liang, Ke Li, Jian Zhang, Meihua Wang, and Chen Lin. Robust visual tracking via identifying multi-scale patches. *Multimed. Tools Appl.*, 78(11):14195–14230, 2019.
56. Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In *Proc. ECCV*, pages 740–755, 2014.
57. Zhiguan Lin and Chun Yuan. Robust visual tracking in low-resolution sequence. In *Proc. ICIP*, pages 4103–4107, 2018.
58. Mingjie Liu, Cheng Bin Jin, Bin Yang, Xuenan Cui, and Hakil Kim. Occlusion-robust object tracking based on the confidence of online selected hierarchical features. *IET Image Proc.*, 12(11):2023–2029, 2018.
59. Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng Yang Fu, and Alexander C. Berg. SSD: Single shot multibox detector. In *Proc. ECCV*, pages 21–37, 2016.
60. Alan Lukežič, Tomáš Vojtíš, Luka Čehovin Zajc, Jiří Matas, and Matej Kristan. Discriminative correlation filter tracker with channel and spatial reliability. *IJCV*, 126(7):671–688, 2018.
61. Chao Ma, Jia Bin Huang, Xiaokang Yang, and Ming Hsuan Yang. Hierarchical convolutional features for visual tracking. In *Proc. IEEE ICCV*, pages 3074–3082, 2015.
62. Chao Ma, Yi Xu, Bingbing Ni, and Xiaokang Yang. When correlation filters meet convolutional neural networks for visual tracking. *IEEE Signal Process. Lett.*, 23(10):1454–1458, 2016.
63. Chao Ma, Jia Bin Huang, Xiaokang Yang, and Ming Hsuan Yang. Robust visual tracking via hierarchical convolutional features. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2018.
64. Chao Ma, Jia Bin Huang, Xiaokang Yang, and Ming Hsuan Yang. Adaptive correlation filters with long-term and short-term memory for object tracking. *IJCV*, 126(8):771–796, 2018.
65. R. J. Mozhdehi and H. Medeiros. Deep convolutional particle filter for visual tracking. In *Proc. IEEE ICIP*, pages 3650–3654, 2017.

66. Reza Jalil Mozhdehi, Yevgeniy Reznichenko, Abubakar Siddique, and Henry Medeiros. Deep convolutional particle filter with adaptive correlation maps for visual tracking. In *Proc. ICIP*, pages 798–802, 2018.
67. Hyeonseob Nam and Bohyung Han. Learning multi-domain convolutional neural networks for visual tracking. In *Proc. IEEE CVPR*, pages 4293–4302, 2016.
68. M.R. Parate, V.R. Satpute, and K.M. Bhurchandi. Global-patch-hybrid template-based arbitrary object tracking with integral channel features. *Appl. Intell.*, 48:300–314, 2018.
69. Shi Pu, Yibing Song, Chao Ma, Honggang Zhang, and Ming Hsuan Yang. Deep attentive tracking via reciprocative learning. In *Proc. NIPS*, pages 1931–1941, 2018.
70. Yuankai Qi, Shengping Zhang, Lei Qin, Hongxun Yao, Qingming Huang, Jongwoo Lim, and Ming Hsuan Yang. Hedged deep tracking. In *Proc. IEEE CVPR*, pages 4303–4311, 2016.
71. Litu Rout, Deepak Mishra, and Rama Krishna Sai Subrahmanyam Gorthi. WAEF: Weighted aggregation with enhancement filter for visual object tracking. In *Proc. ECCVW*, pages 83–99, 2019.
72. Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet large scale visual recognition challenge. *IJCV*, 115(3):211–252, 2015.
73. Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *Proc. ICLR*, pages 1–14, 2014.
74. Yibing Song, Chao Ma, Lijun Gong, Jiawei Zhang, Rynson W.H. Lau, and Ming Hsuan Yang. CREST: Convolutional residual learning for visual tracking. In *Proc. ICCV*, pages 2574–2583, 2017.
75. Chong Sun, Dong Wang, Huchuan Lu, and Ming Yang. Learning spatial-aware regressions for visual tracking. In *Proc. IEEE CVPR*, pages 8962–8970, 2018.
76. Chong Sun, Dong Wang, Huchuan Lu, and Ming Hsuan Yang. Correlation tracking via joint discrimination and reliability learning. In *Proc. IEEE CVPR*, pages 489–497, 2018.
77. Yuxuan Sun, Chong Sun, Dong Wang, You He, and Huchuan Lu. ROI pooled correlation filters for visual tracking. In *Proc. CVPR*, pages 5783–5791, 2019.
78. Fuhui Tang, Xiankai Lu, Xiaoyu Zhang, Shiqiang Hu, and Huanlong Zhang. Deep feature tracking based on interactive multiple model. *Neurocomputing*, 333:29–40, 2019.
79. D.E. Touil, N. Terki, and S. Medouakh. Learning spatially correlation filters based on convolutional features via PSO algorithm and two combined color spaces for visual tracking. *Appl. Intell.*, 48:2837–2846, 2018.
80. Jack Valmadre, Luca Bertinetto, Joao Henriques, Andrea Vedaldi, and Philip H.S. Torr. End-to-end representation learning for correlation filter based tracking. In *Proc. IEEE CVPR*, pages 5000–5008, 2017.

81. Joost Van De Weijer, Cordelia Schmid, and Jakob Verbeek. Learning color names from real-world images. In *Proc. IEEE CVPR*, pages 1–8, 2007.
82. Lijun Wang, Wanli Ouyang, Xiaogang Wang, and Huchuan Lu. Visual tracking with fully convolutional networks. In *Proc. IEEE ICCV*, pages 3119–3127, 2015.
83. Mengmeng Wang, Yong Liu, and Zeyi Huang. Large margin object tracking with circulant feature maps. In *Proc. IEEE CVPR*, pages 4800–4808, 2017.
84. Ning Wang, Wengang Zhou, Qi Tian, Richang Hong, Meng Wang, and Houqiang Li. Multi-cue correlation filters for robust visual tracking. In *Proc. IEEE CVPR*, pages 4844–4853, 2018.
85. Qiang Wang, Jin Gao, Junliang Xing, Mengdan Zhang, and Weiming Hu. DCFNet: Discriminant correlation filters network for visual tracking, 2017. URL <http://arxiv.org/abs/1704.04057>.
86. Xinyu Wang, Hanxi Li, Yi Li, Fatih Porikli, and Mingwen Wang. Deep tracking with objectness. In *Proc. ICIP*, pages 660–664, 2018.
87. Yong Wang, Xinbin Luo, Lu Ding, Jingjing Wu, and Shan Fu. Robust visual tracking via a hybrid correlation filter. *Multimed. Tools Appl.*, 2019.
88. Yi Wu, Jongwoo Lim, and Ming Hsuan Yang. Online object tracking: A benchmark. In *Proc. IEEE CVPR*, pages 2411–2418, 2013.
89. Yi Wu, Jongwoo Lim, and Ming Hsuan Yang. Object tracking benchmark. *IEEE Trans. Pattern Anal. Mach. Intell.*, 37(9):1834–1848, 2015.
90. S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He. Aggregated residual transformations for deep neural networks. In *Proc. IEEE CVPR*, pages 5987–5995, 2017.
91. Yang Yi, Liping Luo, and Zhenxian Zheng. Single online visual object tracking with enhanced tracking and detection learning. *Multimed. Tools Appl.*, 78(9):12333–12351, 2019.
92. Di Yuan, Xinming Zhang, Jiaqi Liu, and Donghao Li. A multiple feature fused model for visual object tracking via correlation filters. *Multimed. Tools Appl.*, 2019.
93. Jianming Zhang, Shugao Ma, and Stan Sclaroff. MEEM: Robust tracking via multiple experts using entropy minimization. In *Proc. ECCV*, pages 188–203, 2014.
94. Peng Zhang, Tao Zhuo, Wei Huang, Kangli Chen, and Mohan Kankanhalli. Online object tracking based on CNN with spatial-temporal saliency guided sampling. *Neurocomputing*, 257:115–127, 2017.
95. Tianzhu Zhang, Changsheng Xu, and Ming Hsuan Yang. Multi-task correlation particle filter for robust object tracking. In *Proc. IEEE CVPR*, pages 4819–4827, 2017.
96. Zhipeng Zhang and Houwen Peng. Deeper and wider Siamese networks for real-time visual tracking, 2019. URL <http://arxiv.org/abs/1901.01660>.
97. S. Zheng, S. Jayasumana, B. Romera-Paredes, V. Vineet, Z. Su, D. Du, C. Huang, and P. H. S. Torr. Conditional random fields as recurrent neural networks. In *Proc. IEEE ICCV*, pages 1529–1537, 2015.

- 
98. Zheng Zhu, Guan Huang, Wei Zou, Dalong Du, and Chang Huang. UCT: Learning unified convolutional networks for real-time visual tracking. In *Proc. ICCVW*, pages 1973–1982, 2018.